



ALGUNOS FACTORES QUE INCIDEN EN EL RENDIMIENTO Y LA EVALUACIÓN EN LOS ALUMNOS DE LAS PRUEBAS DE APTITUD DE ACCESO A LA UNIVERSIDAD (PAAU)

ANNA CUXART JARDÍ (*)
MANUEL MARTÍ RECOBER (**)
FERRAN FERRER JULIÀ (***)

PRESENTACIÓN Y ANTECEDENTES

En junio de 1993 el Boletín Oficial del Estado publicaba la Orden Ministerial de 9 de junio sobre las pruebas de aptitud para el acceso a las facultades, escuelas técnicas superiores y colegios universitarios. En ella se señalaba que: «...una vez finalizado todo el proceso de las pruebas, cuando se observe una elevada desviación entre las medias de los expedientes ¹ de los alumnos y las calificaciones globales otorgadas por un Tribunal,... procederán a estudiar las causas... y a proponer las oportunas medidas en relación con los Centros o Tribunales correspondientes». En el mismo año, la Memoria de Actividades del Consejo de Universidades, insistía en la necesidad de controlar, en relación al centro, la

diferencia entre la media de expediente y la media de las Pruebas de Aptitud de Acceso a la Universidad (PAAU). Ese mismo año, y de manera paralela al trabajo sistemático que se realizaba en la Oficina de Coordinación del COU i les PAAU de Catalunya, Anna Cuxart inició un trabajo de investigación ² dirigido por Manuel Martí Recober. El objetivo consistía en investigar aquellos modelos estadísticos que podían ser de utilidad en el análisis de los resultados académicos de las PAAU, prestando una atención especial al estudio de las variables o factores que inciden en la evaluación de los alumnos que se presentan a dichas pruebas. Las preguntas que entonces se plantearon fueron las siguientes:

- ¿Cuál el grado de asociación entre las puntuaciones que obtienen los

(*) Departament d'Economia, Universitat Pompeu Fabra.

(**) Catedrático de la Universidad Politécnica de Cataluña, Departamento de Estadística e Investigación Operativa. Coordinador del COU y de las Pruebas de Actitud de Acceso a la Universidad (PAAU) en Cataluña.

(***) Departamento de Pedagogía Sistemática y Social, Universidad Autónoma de Barcelona.

(1) La media de Expediente (o *nota de expediente*) de cada alumno es la media aritmética entre los cuatro cursos de secundaria (tres cursos de bachillerato y el COU). La *nota de acceso a la universidad* que servirá más tarde para ordenar los estudiantes y ubicarlos en los estudios superiores es el resultado de calcular la media aritmética entre la nota de expediente y la nota PAAU (media ponderada del conjunto de PAAU).

(2) El resultado de dicha investigación constituye el núcleo de una tesis doctoral en curso, de la cual se han presentado resultados parciales en CUXART, GRAFFELMAN i MARTÍ (1995); CUXART and LONGFORD (1996); MARTÍ y CUXART (1997).

estudiantes (en cada materia y globalmente) en el COU y en las PAAU? Es decir, ¿cuál es la capacidad predictiva de una puntuación respecto de la otra?

- ¿Existe homogeneidad entre los centros que imparten el COU en cuanto a los resultados que obtienen sus alumnos en las PAAU? ¿Y existe esa homogeneidad respecto a la puntuación que estos centros les otorgan? ¿Se puede hablar de puntuación uniforme (estandarizada) en el COU?
- Y en las PAAU, ¿existen diferencias apreciables entre los resultados que dan distintos correctores de una misma materia, hecho que desvirtuaría el principio de prueba estándar?

A continuación damos cuenta de algunas de las características que tuvo la investigación que se llevó a cabo para dar respuesta a las anteriores preguntas:

- La modelización estadística, construida a partir de datos individuales.
- La utilización de amplia información ³ sobre cada alumno: notas del COU y de las PAAU por asignaturas —ésta es una de las diferencias a destacar respecto de otros estudios publicados—, características personales del alumno y del centro en el que éste estudió el COU, etc.
- Una investigación empírica llevada a cabo para detectar y cuantificar las imperfecciones del proceso.
- La propuesta de ajustes incorporados al proceso (autorevisión) con el objetivo de obtener una estimación más eficiente de la preparación de los alumnos. En este sentido se se-

ñalan ciertas perspectivas de automatización en el futuro.

Con este trabajo pretendemos reflexionar sobre los factores que inciden en el rendimiento y en la evaluación de los alumnos en las PAAU. En el primer apartado se estudia la influencia del centro escolar en la predicción de la nota PAAU. En el segundo apartado se analiza la influencia del profesor-corrector en el proceso. Y por último, en el espacio dedicado a las conclusiones, se recogen algunas perspectivas de investigación y otras consideraciones pedagógicas que han surgido a partir de los resultados estadísticos.

LA INFLUENCIA DEL CENTRO ESCOLAR EN LA PREDICCIÓN DE LA NOTA PAAU DE CADA ALUMNO

ANTERIORES ESTUDIOS

Muchos han sido los investigadores que anteriormente se han dedicado a estudiar en profundidad las llamadas pruebas de Selectividad. Entre los pioneros, destacan T. Escudero e I. Aguirre de Cárcer. En el documento del Consejo de Universidades (que antes ha sido citado) se recoge una amplia bibliografía que permite conocer cuáles han sido las preocupaciones de los investigadores sobre el tema, sus principales aportaciones y las sucesivas modificaciones que se han ido introduciendo en la normativa de las pruebas a la luz de dichas investigaciones. Aquí nos limitaremos a mencionar aquellos trabajos que por sus objetivos y resultados están directamente relacionados con los estudios que presentamos.

(3) La Oficina de Coordinació del COU i les PAAU de Catalunya dispone de las notas por asignaturas de todos los alumnos que han cursado COU. Puesto que el examen PAAU se basa en las materias cursadas en COU nos ha parecido más adecuado analizar la relación entre estas puntuaciones que entre la nota de expediente y la de las PAAU.

Antoni Sans (1990) analizó los resultados en las PAAU que obtuvieron los 12.423 estudiantes que se matricularon en el COU en el curso 1986-87 y que estaban adscritos a la Universidad Autónoma de Barcelona. Su estudio señala diferencias significativas en la nota PAAU entre los distintos tribunales y las distintas convocatorias de exámenes (tandas). Sans llega a proponer que no se tenga en cuenta la nota de expediente en el cómputo de la nota de acceso a la universidad a la luz de las discrepancias observadas entre nota de expediente y nota PAAU en algunos centros. Precisamente los centros con peores resultados en las PAAU eran los que habían puntuado más alto a sus alumnos. El manifiesto comportamiento heterogéneo que Sans observa entre los centros le lleva a esta recomendación. El desequilibrio existente se manifiesta también entre centros privados y públicos: mientras que el porcentaje de los alumnos matriculados en los centros públicos es de un 50 *por ciento* para los que superan el COU, y de un 38 *por ciento* para los que aprueban las PAAU, en los centros privados estos porcentajes ascienden a un 78 *por ciento* y a un 66 *por ciento*, respectivamente.

Otro trabajo a señalar es el de Mercedes Muñoz y otros autores (1991) realizado a partir de los principales resultados de sus investigaciones sobre las pruebas de acceso a la universidad. En dicho trabajo, también

estudian la incidencia de las modificaciones que se introdujeron en el formato y organización de las PAAU a partir de los resultados de las convocatorias de 1987, 1988 y 1989 relativos a una muestra estratificada de universidades de todo el Estado. Nos parece interesante destacar aquí alguna de las conclusiones de este trabajo en cuanto al sistema actual de acceso a la universidad: «...Podrían resumirse sus defectos en que cumple mal y de forma desigual la función que se le asigna de ordenar de manera aquilatada ⁴ a los estudiantes en cuanto a su prioridad para obtener un puesto en la universidad».

LOS DATOS. UNA PRIMERA EXPLORACIÓN

Para abordar el estudio de la asociación de la nota que obtuvo cada estudiante en las PAAU y en el COU ⁵, se escogió una muestra aleatoria ⁶ de 26 centros (1.619 estudiantes) del distrito de Catalunya, sobre los 400 centros (unos 25.000 estudiantes en total) que concurrieron a las PAAU en junio de 1993. Nuestra hipótesis de trabajo era que la relación entre la nota obtenida en el COU y la de las PAAU variaba de un centro a otro. Se trataba pues de escoger un modelo de variación de la nota PAAU que incluyera la nota COU como variable explicativa (en-

(4) La preocupación sobre el grado de homogeneidad existente entre los centros al puntuar a sus alumnos sigue siendo motivo de preocupación en uno de los estudios publicados recientemente. M.^a del R. LÓPEZ (1997) estudia los resultados de las pruebas de acceso a la Universidad Autónoma de Madrid del curso 1995-1996. La distribución por centros (132 centros de secundaria) de la diferencia entre la nota de expediente y la nota de las PAAU presenta un promedio de -1.6 con una variabilidad que va desde -0.8 hasta -2.8. Para los estudiantes del bachillerato LOGSE el promedio de dicha diferencia por centros (14 centros) es de -1.9 variando entre -0.9 y -3. Luego, no parece que el nuevo sistema educativo vaya a reducir estas diferencias en la evaluación de los alumnos. Sigue pues siendo un tema pendiente.

(5) Nota COU es la media aritmética de las ocho asignaturas pertenecientes a este curso mientras que nota PAAU, como decíamos antes, es una media ponderada de los nueve ejercicios que componen dichas pruebas en Catalunya.

(6) La Oficina de Coordinació del COU i les PAAU de Catalunya suministró los datos.

tre otras) y que a la vez contemplara la posible diferencia entre centros.

En una primera exploración constatamos que la correlación intra-centros ⁷ para la nota PAAU era de 0.195. Es decir, aproximadamente, un 20% de la variación de la nota PAAU observada se debía a variación entre centros. Este hecho desaconsejaba la aproximación clásica, a través de un modelo de regresión con una misma ecuación para los 1.619 estudiantes ⁸.

En el Gráfico I se observa:

- La devaluación de nota (del COU a las PAAU) de la mayoría de los estudiantes. Si la nota del COU y la nota de las PAAU estuvieran midiendo la misma variable latente «aptitud del alumno para los estudios universitarios», sería de esperar que el gráfico de dispersión se situara alrededor de la recta bisectriz de este primer cuadrante. Pero no es así, los puntos se sitúan alrededor de una recta paralela a dicha bisectriz. Existe, por tanto, un sesgo entre una valoración y la otra. Está claro que esta diferencia se distribuye de manera desigual entre los estudiantes. Pero nos

preguntamos también, ¿cómo se distribuye entre centros?

- Las distribuciones marginales de la nota obtenida en el COU y de la nota de las PAAU. Más de la mitad de los estudiantes tienen una nota COU comprendida entre los 6 y 7 puntos. Cabe recordar que los alumnos que se presentan a las PAAU deben haber aprobado el COU, es decir han de tener una nota COU superior a un 5.5. Por otro lado, un 65% de los estudiantes obtiene en las PAAU una nota situada entre los 4 y 6 puntos; un 8.4% obtiene una nota inferior a 4; un 20.7% consigue una puntuación entre 6 y 7; un 6.2% entre 7 y 9; y un 0.1% obtiene una calificación superior a 9.

Si en el Gráfico I nos centramos solamente en los resultados obtenidos por los estudiantes de dos escuelas concretas (la número 2 y la número 18 de nuestra muestra), observamos la situación que se aprecia en el Gráfico II.

El Gráfico II sugiere la conveniencia de estimar un modelo de regresión con parámetros que puedan tomar valores diferentes ⁹ para cada escuela.

(7) El coeficiente de correlación intra-grupos es una medida del grado de homogeneidad existente dentro de los grupos. La estimación que hemos tomado es la que propone MUTHÉN (1994) a partir de la descomposición derivada de un análisis de varianza clásico. Para más detalle ver COCHRAN (1977).

(8) Al no satisfacer la hipótesis de incorrelación entre residuos, los estimadores de cuadrados mínimos (MCO) subestiman los errores estándar de los coeficientes de la ecuación de regresión. Esta «subestimación» de los errores estándar conlleva que no puedan darse por válidos los test que suelen ofrecer los paquetes de software estadístico al realizar un análisis de la regresión. Al mismo tiempo, y puesto que nuestro interés se centraba en toda la población de centros, no únicamente en los 26 de la muestra, desestimamos los modelos de análisis de la covarianza (ANCOVA) que ofrecen la estimación del efecto fijo debido a cada centro sin inferir sobre la población de los mismos.

(9) Si aplicáramos un modelo de regresión clásico por separado a estas dos escuelas se obtendrían dos rectas diferentes quizás paralelas, es decir, con la misma pendiente y diferente ordenada en el origen.

GRÁFICO I
Nota PAAU versus nota COU (1619 estudiantes)

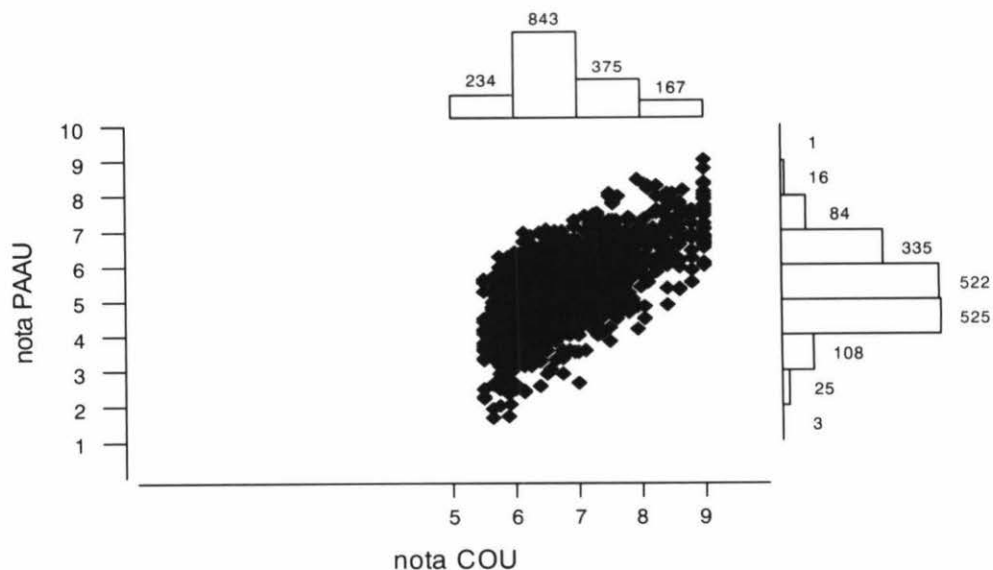
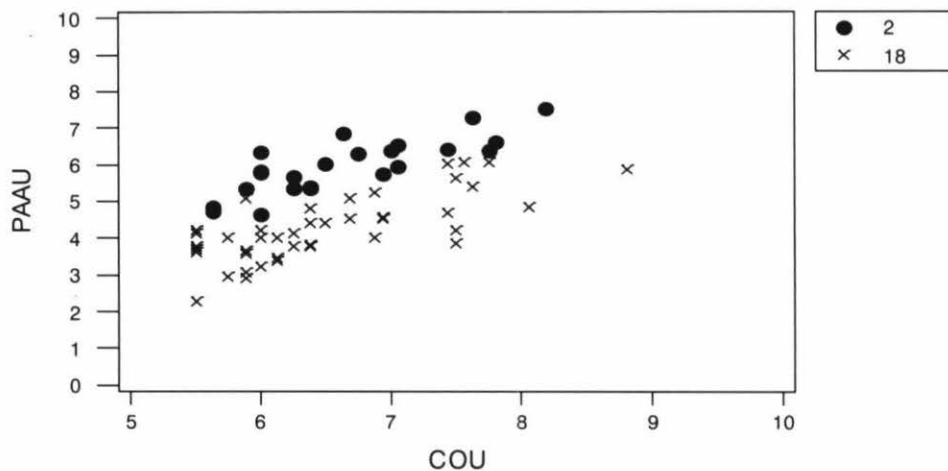


GRÁFICO II
Notas de los alumnos de los centros 2 y 18. Las notas medias del COU y de las PAAU son, respectivamente: 6.62 y 5.91, el primero; 6.46 y 4.17, el segundo.

Nota global PAAU versus nota global COU (escuelas 2 y 18)



ASOCIACIÓN ENTRE LA NOTA DEL COU Y LA NOTA DE LAS PAAU

Con el objetivo de poder hacer inferencias sobre la población de centros y a la vez estimar eficientemente los parámetros de interés, decidimos *modelizar* la variación de la nota PAAU a través de un sistema de regresión de coeficientes aleatorios (parámetros que varían de una escuela a otra de acuerdo con una distribución de probabilidad) tomando la nota del COU como variable explicativa. Aitkin and Longford (1986) proponen la utilización de dichos modelos para el estudio de datos de estructura jerárquica. En nuestro caso, los estudiantes (o unidades del primer nivel) aparecen agrupados en centros escolares (o unidades de segundo nivel). El método de estimación seguido ha sido el de Mínimos Cuadrados Generalizados It-

rativos (MCGI) y se ha utilizado el software informático Mln (ver Goldstein, 1995).

En la tabla I presentamos los resultados que se derivan del ajuste de diferentes modelos a los datos de la muestra así como las variables incorporadas en cada uno de estos modelos. El modelo 1 es un análisis de la regresión clásico en el que no se tiene en cuenta la agrupación de los estudiantes en centros. El modelo 2 es el más simple de los modelos de coeficientes aleatorios de dos niveles: descompone la varianza observada de la nota PAAU en varianza entre centros y en varianza entre alumnos dentro de los centros; estima la media global de la nota obtenida en las PAAU (5.24) y permite además calcular una estimación del coeficiente de correlación intra-centros. Por su parte, el modelo 3, amplía el modelo anterior incluyendo la variable explicativa nota COU. Como re-

TABLA I

La variación de la nota PAAU individual con relación a la nota COU, el género, el ser alumno repetidor de COU y la opción

	Modelo 1 MCO	Modelo 2 MCGI	Modelo 3 MCGI	Modelo 4 MCGI
<i>Coefficientes (E.E.)</i>				
constante	-0.84 (0.17)	5.24 (0.10)	-0.68 (0.17)	-0.36 (0.18)
NOTCOU	0.90 (0.03)		0.88 (0.02)	0.85 (0.02)
SEXO				0.20 (0.02)
REPCOU				-0.23 (0.05)
OPA				-0.27 (0.04)
OPB				-0.41 (0.05)
<i>Varianzas</i>				
entre centros	—	0.22 (0.07)	0.18 (0.07)	0.18 (0.05)
entre estudiantes	0.706	1.03 (0.04)	0.52 (0.02)	0.48 (0.02)
<i>Coef. de correlación</i>				
<i>intra-centros</i>	—	0.175	0.25	0.27

sultado de la estimación ¹⁰ podemos decir que el coeficiente de dicha variable (pendiente de la recta) es el mismo para todos los centros, no varía. Sin embargo, la constante (ordenada en el origen) varía entre los distintos centros. Esta constante, que varía de un centro a otro, está recogiendo el efecto producido debido al centro escolar ¹¹. El modelo 3 se amplió añadiendo diversas variables cualitativas tales como el género, el ser repetidor de COU, la opción de COU, el tipo de centro (público o privado). El tipo de centro no resultó ser una variable significativa. No se observa diferencia entre la opción C y la D, que tomamos como opción base o referencia. Por otra parte, el modelo más completo y que, además de la variable nota COU, incorpora otras variables que explican la variación de la nota PAAU es el Modelo 4. La exploración de los residuos del modelo 4 puso en evidencia ninguna violación de las hipótesis sobre los mismos con lo cual es un buen modelo para explicar la variación de la nota PAAU. Según estos dos modelos (3 y 4), el efecto centro en la nota PAAU es significativo y se reduce a un término aditivo de media 0 y varianza 0.18. En consecuencia, para cada centro la recta de regresión de PAAU sobre COU tiene una constante propia, resultante de sumar la constante común a todos los centros (-0.68 en el modelo 3) con el efecto centro correspondiente, y una pendiente fija (0.88 en el modelo 3), la misma para todos los centros. Se confirma, pues, el patrón de comportamiento que sugería el Gráfico II.

Es de destacar que, al ir añadiendo variables explicativas (modelos 2, 3 y 4) se va

reduciendo la varianza entre centros y entre estudiantes -esta última en mayor grado-, y que sin embargo, el coeficiente de correlación intra-centros aumenta. Al ajustar (corregir) la variación de la nota PAAU con las sucesivas covariantes se hace todavía más ostensible la homogeneidad interna de los centros frente a la heterogeneidad de los mismos en lo que se refiere a la puntuación en las PAAU.

Centrándonos en la interpretación del modelo 4, se puede observar que la nota obtenida en el COU es predictora ¹² de la puntuación que se puede conseguir en las PAAU (matizada con una serie de características individuales); pero también se observa que el factor centro escolar tiene carácter predictivo. Por otro lado, es importante destacar ciertas observaciones al respecto:

- Por el hecho de haber cursado el COU en un centro u otro y respecto al comportamiento medio de la población, la nota PAAU que se espera de cada estudiante sufrirá un incremento que puede oscilar entre -0.84 y 0.84 puntos ¹³.
- Se observa que en las PAAU, en igualdad de condiciones con respecto al resto de variables, a los chicos les va mejor (0.20 puntos en promedio) que a las chicas.
- El hecho de ser repetidor disminuye la predicción de nota en 0.23 puntos.
- A los estudiantes de Ciencias (0.27 por debajo de la media para los de la opción A y 0.41 para los de la opción B) les va peor que a sus compañeros de Letras (opciones C y D).

(10) Hemos desestimado el modelo de dos niveles que contempla la ordenada en el origen y la pendiente variando de un centro a otro ya que la varianza de este último coeficiente no resultó significativa (no se puede rechazar la hipótesis de que dicho coeficiente sea constante) y además dicho modelo ofrecía un peor ajuste (test de la razón de verosimilitud) en comparación con el Modelo 3.

(11) A partir de este momento, para abreviar, nos referiremos al efecto debido al centro estimado en los modelos 2, 3 y 4 como el *efecto centro*.

(12) El coeficiente de correlación lineal entre la nota COU y la nota PAAU de cada estudiante es 0.66.

(13) $0.84 = 2DE = 2\sqrt{0.18}$

EL EFECTO DEBIDO AL CENTRO ESCOLAR Y LAS ESCALAS DE MEDIDA

Los resultados de la estimación del modelo 4 corroboran los que, para el ámbito nacional, encontraron Muñoz-Repiso (1991) y otros autores. Nuestro enfoque, realizado a través de modelización estadística, ofrece, además, la posibilidad de obtener un identificador para cada centro¹⁴. Los centros que presentan una asociación extrema entre la nota obtenida en las PAAU y la nota del COU deberían ser motivo de estudio. El análisis de las características de dichos centros y la discusión con los responsables de los mismos podría aportar un mayor conocimiento sobre la diversidad¹⁵ de los centros.

Mientras la nota PAAU presenta una variación considerable entre distintos centros, no ocurre lo mismo con la nota del COU. El coeficiente de correlación intra-centros para esta última nota es prácticamente 0. Se podría decir que las distribuciones de puntuaciones (aprobados) en el COU no varían de un centro a otro. Cada centro ha «ordenado» a sus alumnos con una media y una variabilidad muy similares. En cambio, según acabamos de ver, al proponer el mismo examen

a todos los alumnos se aprecian las diferencias entre los resultados de los alumnos de un centro a otro. Nuestra conclusión es que los centros están utilizando escalas de medida diferentes y distintos estándares en la preparación de los alumnos.

LA INFORMACIÓN A LOS CENTROS. INDICADORES E INFORMACIÓN

Los estudiantes concurren a las PAAU a través del centro en que han cursado el COU. Cada centro está adscrito a una determinada universidad. Por ejemplo, en Cataluña, la Oficina de Coordinació del COU i les PAAU es la encargada de coordinar y administrar la realización de estas pruebas (los exámenes son los mismos para todas las universidades) y, a la vez, tiene como función la de hacer llegar a los centros los resultados de sus alumnos. La información que suministra se concreta en la nota media del COU, la nota media obtenida en las PAAU y la diferencia entre ambas notas para los alumnos de su centro. Además comunica la nota media del conjunto de alumnos de la convocatoria. En la actualidad

(14) Dicho identificador, al tener asociada, según el modelo, una distribución de probabilidad, permite distinguir los valores extremos de los considerados más comunes. En los modelos de regresión de nivel múltiple (GOLDSTEIN, 1995) nos encontramos con residuos asociados a cada una de las variables o coeficientes aleatorios introducidos en el modelo. El modelo que hemos estimado para nuestros datos contempla la existencia de residuos a nivel individuo (tantos como individuos) y residuos a nivel centro (tantos como centros). Cada residuo no es más que una realización (o predicción de realización) de una variable aleatoria. El residuo componente a cada centro (*efecto centro*) se estima a partir de los residuos individuales, el número de alumnos y el valor del coeficiente de correlación intra-centros. Se trata de *estimadores encogidos* (subestiman la información sobre un subgrupo a favor de la información sobre todo el grupo cuando el número de elementos del subgrupo es muy pequeño y puede conducir a estimaciones poco eficientes). Una vez identificado cada centro con su *efecto* podemos establecer una banda de efectos no distinguibles (que se separan como máximo una desviación estándar de la media) y considerar aquellos centros que presentan un efecto superior a 0.42 o inferior a -0.42.

(15) Los centros que impartirán el bachillerato LOGSE proceden de institutos de bachillerato, institutos de formación profesional o centros de enseñanza secundaria creados *ex profeso*. Cabe esperar, pues, diferencias importantes entre los centros, no solamente en cuanto al conjunto de materias sino también en cuanto al estilo de la enseñanza que impartirán.

éstos son algunos de los indicadores ¹⁶ con los que cuenta la comunidad educativa y cómo tales indicadores han sido utilizados desde la Administración y desde cada uno de los centros para realizar comparaciones.

Nosotros, nos preguntamos ahora sobre la «oportunidad» de tales comparaciones ¹⁷ dado que, en muchos casos, las diferencias observadas no parecen ser significativas. Es más, debido a la diversidad existente en cuanto al número de alumnos ¹⁸ que cada centro presenta a las PAAU, algunas comparaciones carecen incluso de sentido.

Por estos motivos, a la hora de abordar esta cuestión, decidimos estudiar la estructura de covarianza entre la nota COU y la nota PAAU descomponiendo la variación global en variación entre centros y variación dentro de los centros. La metodología seguida consiste en *modelizar* de manera conjunta (bivariante) la variación de la nota COU y la nota PAAU de cada estudiante, descomponiendo la variación de cada nota en efecto debido al centro y efecto específico del estudiante y admitiendo la posible covarianza entre los efectos de un mismo nivel. Hemos aplicado a los datos de la muestra (que se recogían en el apartado que llevaba por título *Los datos. Una primera aproximación*) un

modelo de componentes de la varianza para el estudio de datos multivariantes con estructura jerárquica ¹⁹. Los resultados ²⁰ de la estimación de la variación entre centros, matriz Σ_c , y dentro de los centros entre estudiantes, matriz Σ_E son los siguientes:

$$\Sigma_c = \begin{pmatrix} 0.022 & 0 \\ 0 & 0.201 \end{pmatrix} \quad \Sigma_E = \begin{pmatrix} 0.662 & 0.582 \\ 0.582 & 1.033 \end{pmatrix}$$

CONCLUSIONES

Una primera conclusión que se desprende del tratamiento estadístico señala que la media global de la nota COU y de la nota PAAU es, respectivamente, 6.74 y 5.24. La descomposición de la variación total en variación entre centros y dentro de los centros está justificada. La varianza entre centros es significativa para ambas notas. Los resultados de estas estimaciones corroboran lo que desde el punto de vista descriptivo ya se apuntaba anteriormente: una mayor variabilidad entre efectos centro en la nota PAAU que en la nota COU. Mientras el efecto centro en la nota PAAU toma valores en un rango que va de -0.9 a 0.9 ²¹, en la nota COU varía entre -0.3 y 0.3. En el Gráfico III puede observarse que el

(16) En este sentido, hay que celebrar la labor que está realizando el INCE en su «Proyecto de sistema estatal de indicadores de la educación». Otra referencia imprescindible es el documento de la OCDE «Education at a glance. OCDE indicators» de 1995.

(17) Estos datos no se hacen públicos en Cataluña –tampoco nos consta que se haga en el resto de España– hasta el momento. En otros países, como Inglaterra y el País de Gales donde tradicionalmente se publican estas informaciones y, a la vista del uso no siempre adecuado que se ha hecho de estos indicadores, los científicos (ver artículo de GOLDSTEIN y SPIEGELHALTER, 1996) se han visto obligados a recordar la incertidumbre inherente en este tipo de medidas y la necesidad de contextualizarlas.

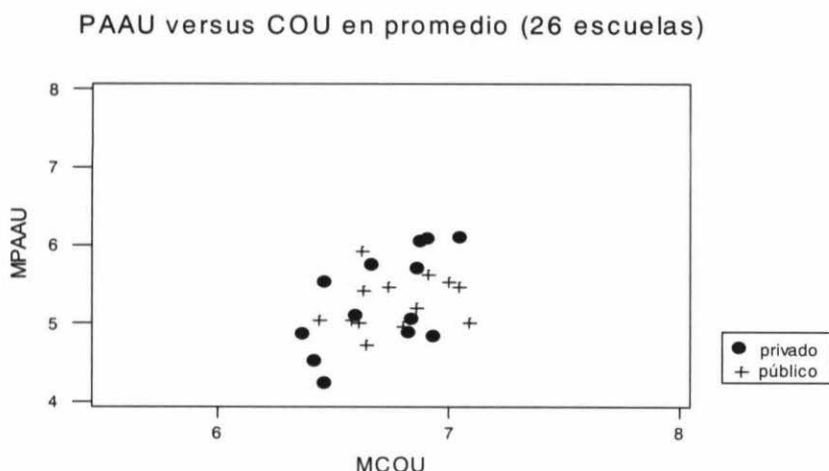
(18) Los centros con un número reducido de estudiantes son más susceptibles de verse afectados en sus medias por valores excepcionales. El tamaño del centro es una característica que se debería tener en cuenta al comparar centros.

(19) Para más detalle sobre este tipo de modelos ver GOLDSTEIN (1987).

(20) La diferencia con el modelo del apartado anterior radica en que, en aquel caso, tomábamos una única respuesta, la nota PAAU y estudiábamos su variación previo ajuste con la nota COU. Ahora consideramos ambas notas como respuestas y estudiamos su variación conjunta. Los detalles sobre este enfoque bivariante de nivel múltiple y su aplicación a la muestra de 26 centros se encuentran en la tesis que A. Cuxart está elaborando.

(21) $0.90 = 2DE = 2\sqrt{0.201}$

GRÁFICO III
Medias en COU y en las PAAU de los 26 centros



rango de variación de la nota media de las PAAU es mucho mayor (más del doble) que el de la nota media del COU. Tanto desde la consideración del alumno, como desde la del centro, la nota PAAU presenta una mayor variabilidad ²² que la nota COU. No solamente la primera nota distingue mejor los centros, sino que, dentro de un mismo centro, la nota PAAU introduce más discriminación que la nota COU.

Una segunda conclusión igualmente importante resulta de este aspecto: atendiendo al nivel del centro, la covarianza no es significativa. La covarianza nula entre el efecto centro en la nota COU y el correspondiente en la nota PAAU nos indicaría que el hecho de tener una media alta de la nota COU no siempre va acompañado de una media alta también en la nota PAAU. En cambio, en relación al alumno, la correlación entre efectos es

0.704 (superior a 0.66, la correlación obtenida entre las dos notas a nivel global). Existe mayor coherencia entre ambas puntuaciones dentro de cada centro que si consideramos a todos los alumnos de la muestra. Este hecho tiene una clara explicación: cada centro ha ordenado a los estudiantes que superan el COU a través de sus puntuaciones y el resultado son distribuciones de notas COU con una media similar de un centro escolar a otro. En las PAAU se realiza una nueva ordenación. Aunque de manera interna en cada centro exista una coherencia entre ambas ordenaciones, parte de la misma se pierde al agrupar a todos los centros ya que la escala o baremo utilizado por cada uno de los centros en COU no es exactamente la misma.

Por último, hay que destacar hasta qué punto resulta fiable la comparación de los distintos centros a través de sus respectivas

(22) La similitud de medias de COU entre centros y dentro de los centros es una consecuencia, en parte, del reducido rango de valores con que se puntúa cada asignatura: 5.5, 6.5, 7.5 y 9. Como decíamos en el apartado anterior, la nota PAAU evidencia que los centros están puntuando con escalas diferentes en COU pero dando como resultado distribuciones de notas similares. En el bachillerato LOGSE se ha corregido esta deficiencia y la puntuación es más «fina».

notas media de COU y de las PAAU. A este respecto y a tenor de nuestros resultados podemos señalar que:

- No tiene sentido comparar los distintos centros en cuanto a las notas medias obtenidas en COU ²³ (en el caso de que dos centros se diferenciara en una cantidad apreciable, esta diferencia solamente sería observable si el tamaño de los centros fuera superior a 197).
- Para centros de tamaño inferior a 30 alumnos tampoco se recomienda establecer ordenaciones a partir de las medias obtenidas en las PAAU.
- En cambio, la ordenación «más informativa», la que puede ser utilizada para tamaños incluso de 16 alumnos, es la diferencia entre las dos medias. Esta ordenación es la que puede resultar más útil para la Inspección Educativa en el estudio de los casos que se separan considerablemente del comportamiento estimado como promedio.

LA INFLUENCIA DEL PROFESOR-CORRECTOR. ANÁLISIS DEL PROCESO DE CORRECCIÓN DE LAS PRUEBAS PAAU ²⁴

ESTUDIOS ANTERIORES Y MOTIVACIÓN

En educación —en general en todo el ámbito de las ciencias sociales— las variables suelen ser difíciles de medir. Por ejemplo, si pretendemos evaluar la habilidad que una población de estudiantes manifiesta en matemáticas, tendremos que defi-

nir previamente aspectos como éstos: qué entendemos por habilidad en matemáticas, qué tipo de prueba prepararemos para provocar que se manifieste tal habilidad (si se trata de un examen, si van a ser preguntas abiertas o cerradas, si se trata de una prueba oral o escrita, etc.), qué respuestas esperamos obtener y cuáles de ellas daremos por válidas, cómo se administrará la prueba, cómo se puntuará, y por fin, una vez tengamos la puntuación final, habrá que definir qué interpretación corresponderá a los posibles resultados. Cada uno de los elementos que integran el proceso de medida conlleva arbitrariedad e incertidumbre. Ante esta situación, digamos de imperfección del instrumento de medida, y puesto que erradicarla es imposible, lo aconsejable es avanzar en el conocimiento de sus causas para limitar al máximo su impacto.

Sans (1991) entre otras consideraciones apunta la necesidad de medir la fiabilidad del proceso de corrección de las PAAU a la vista de las diferencias observadas entre distintos tribunales. Por su parte, Muñoz-Repiso y otros autores (1991) también abordan el tema de la corrección. Así, por ejemplo, al intentar explicar por qué los estudiantes de Ciencias (opciones A y B) sufren una disminución de nota —de la nota de expediente a la nota PAAU— superior a sus compañeros de Letras (opciones C y D) sugieren que una de las razones podría ser que las asignaturas específicas del área de Ciencias permiten una mayor discriminación entre los alumnos que las asignaturas específicas del área de Letras, a la vez que, en las asignaturas comunes a todas las opciones —como Filosofía o Comentario de Texto—, se tiende a discriminar poco otorgando puntuaciones que en su

(23) En el Gráfico III observamos la similitud en las medias de COU para los centros. Llama la atención que los centros que se destacan por obtener los mejores o los peores resultados en media en las PAAU son todos privados. Estos resultados coinciden con los obtenidos por MUÑOZ-REPISO y otros autores (1997) en un estudio relativo a los centros de la Universidad Autónoma de Madrid.

(24) Este apartado ha sido redactado partiendo de una investigación ya presentada que dirige el Sr. Manuel Martí Recober y que está financiada en parte por el CIDE (Proyectos de Investigación, 1995).

mayoría van del 4 al 7. Por ello, y ante la instauración del tribunal único para el curso 1991-1992, Muñoz-Repiso y otros autores (1991) planteaban la necesidad de criterios de homologación y garantías de mayor objetividad en la corrección de las pruebas.

Para medir la fiabilidad de la corrección es necesario disponer de algún tipo de réplicas. Escudero y Bueno (1994) realizaron un experimento con un tribunal paralelo que evaluaba los exámenes puntuados a su vez por el tribunal oficial correspondiente. Al comparar las puntuaciones que otorgaron los dos tribunales, no encontraron diferencias significativas entre las medias de la nota final de las PAAU (tampoco resultaban significativas las diferencias entre la mayor parte de notas agregadas)²⁵. Desde nuestro punto de vista uno de los resultados más relevantes del trabajo de Escudero y Bueno (1994) consistió en la constatación del siguiente hecho: si se hubieran tomado las puntuaciones de la segunda corrección (experimental) en lugar de las oficiales, para un 10% de los estudiantes se hubiera invertido la resolución (obtener o no un aprobado) de la nota de acceso a la universidad.

Aunque de manera global los resultados de un tribunal no sean significativamente diferentes de los de otro tribunal, ésta, creemos, no es una razón suficiente para pensar que el proceso en sí mismo es «justo». Deberían ser los expertos los que se pronunciaran al respecto. Y a partir de

sus observaciones, sería la política educativa la que tendría que definir qué es lo que se entiende por una «diferencia aceptable». Además, los datos de interés en educación deberían considerarse en relación al individuo, y no solamente en cuanto a la agrupación de individuos, ya se trate de un tribunal o de centro escolar.

En junio de 1991, Albert Satorra y Frederic Udina de la UPF, llevaron a cabo un experimento²⁶ de control de la calidad en la corrección de los exámenes de Matemáticas I de las PAAU. De los resultados se pudo estimar que la varianza inducida por la corrección en la nota de Matemáticas I era del 10%. No obstante, insistían sus autores que este estudio no podía ser considerado como concluyente, sino más bien como una invitación a la reflexión, dado el carácter voluntario de las respuestas.

Así se perfilaba la necesidad de abrir una línea de investigación en nuestro país sobre un tema que, hasta ese momento había sido insuficientemente estudiado:²⁷ la necesidad, por una parte de medir eficientemente la variabilidad de la corrección en cada una de las pruebas y por otra, de indagar sobre las componentes de dicha variabilidad.

El trabajo de Albert Satorra y Frederic Udina había constatado la existencia de un efecto debido al centro —centro en que se cursa el COU— en la nota PAAU de cada estudiante. El efecto debido al centro también aparecía al analizar la asociación entre la nota PAAU de cada materia y la co-

(25) En el trabajo de Escudero y Bueno no se comparan los resultados de la doble corrección en las pruebas específicas por asignaturas. Se estudian los resultados a nivel global de la primera obligatoria, segunda obligatoria, primera optativa y segunda optativa.

(26) A los 73 correctores de dicha asignatura de la ciudad de Barcelona se enviaron por correo dos exámenes fotocopiados (al azar se escogieron 20 exámenes de uno de los tribunales y se fotocopiaron antes de ser corregidos oficialmente) pidiéndoles que los corrigieran con el mismo criterio que días antes habían aplicado en la corrección oficial. De los 73 correctores, respondieron 39.

(27) Lo cierto es que a pesar del valioso estudio realizado por el Consejo de Universidades en 1993 y de las interesantes recomendaciones que en el mismo ya se incluían, poco se ha avanzado en estos cuatro años para implementarlas. Una de las recomendaciones era la de utilizar un mayor número de preguntas con la intención de abarcar mejor el programa y así incrementar la fiabilidad de las pruebas.

respondiente nota obtenida en COU. El paso siguiente de nuestra investigación consistió en analizar ²⁸ las posibles causas de dicho efecto. Había que tener en cuenta aspectos, como por ejemplo, el hecho de que, debido a cuestiones organizativas, en Cataluña es el mismo corrector el que corrige los exámenes de todos los alumnos de un mismo centro. Por ello, el efecto centro podía estar confundido con el efecto debido al corrector. Aplicando modelos de regresión de nivel múltiple pudimos separar (Cuxart, 1995) el efecto centro del efecto corrector en la materia de Matemáticas I, estimando una varianza significativa entre centros del mismo orden que la varianza entre correctores ²⁹. Cabe destacar que la magnitud de la incidencia del efecto centro en la nota PAAU resultó ser muy próxima a la obtenida por Satorra y Udina (1991). Pero, el hecho de que cada examen fuera corregido por una sola persona no nos permitía investigar sobre los posibles grados de severidad de los correctores ni tampoco abordar temas de fiabilidad en la corrección.

Con la intención de poder estudiar con más profundidad el proceso de corrección en la calificación de los estudiantes que se presentan a las PAAU, emprendimos el diseño de un experimento que permitiera evaluar la calidad de la corrección en dos materias, en Matemáticas I y en Filosofía (consideradas de diferente dificultad en la concreción y aplicación de los criterios de

corrección). El objetivo principal del estudio consistía en la obtención de medidas de la calidad de la corrección que nos permitieran iniciar un proceso de seguimiento y control del sistema en posteriores convocatorias. Un segundo objetivo que perseguíamos, ligado al anterior y que adquiere sentido en función de él, era la detección de las posibles fuentes o factores de la variabilidad de la corrección.

En evaluación educativa en general y en las pruebas PAAU en particular, interesa que el examen sea válido ³⁰, es decir que el examen mida aquello que ha de medir y para lo que ha sido concebido y que la puntuación que otorga el proceso de corrección sea fiable. La fiabilidad tiene sentido en un contexto de réplicas y se refiere a la precisión del instrumento de medida. En consecuencia, la puntuación será más fiable o precisa cuanto menor sea el error de medida introducido en el proceso de corrección.

En el estudio de Satorra y Udina (1994) la selección de correctores no fue aleatoria. La corrección se llevó a cabo con posterioridad a las pruebas y sin la presión del volumen de exámenes a corregir. El modelo de componentes de la varianza de estos autores no distinguía entre posibles fuentes de error en la estimación. El trabajo de Escudero y Bueno (1994), en el que por el contrario se respetaron las condiciones de realización que acabamos de citar, involucra pocos correctores. Cada tribunal, en aquellos momentos, solía tener so-

(28) Es evidente que para estudiar el efecto centro se debería investigar en el proceso de evaluación en las escuelas. Este sería tema de otra investigación muy interesante, por cierto. Nuestros esfuerzos se centraron en conocer el funcionamiento del instrumento de medida de la nota PAAU, es decir, el proceso de corrección.

(29) Los detalles de este trabajo se encuentran en un informe preparatorio.

(30) No abordamos el tema de la validación del examen. Nos centraremos en el proceso de corrección y en su fiabilidad. No obstante, como se verá más adelante, al analizar las causas de las discrepancias observadas entre correctores, se apunta la posibilidad de que los exámenes propuestos no estén midiendo adecuadamente la preparación de los alumnos. De ser así, tendríamos que una insuficiente validación del examen conllevaría al mismo tiempo una fuente de discrepancia en la corrección, añadiendo más elementos de injusticia al proceso.

lamente un corrector para cada asignatura. Al comparar dos tribunales, en realidad, lo que se estaba haciendo era comparar dos correctores.

Frente a estos estudios, nuestro diseño reunía las siguientes características:

- Se hizo una corrección doble. Cada examen debía ser corregido dos veces, una sería la corrección oficial y la otra la realizaría un corrector adscrito a otro tribunal de las PAAU. Ésta, llamémosle, segunda corrección –aunque se hiciera al mismo tiempo y sin conocimiento de la corrección oficial– se haría a partir de una fotocopia ³¹.
- Se asignó el segundo corrector al azar, no de manera voluntaria.
- Los exámenes fueron corregidos en las fechas oficiales, no con posterioridad.
- Participaron un número importante de estudiantes (187 en Matemáticas I y 363 en Filosofía) lo cual permitió contar con un abanico de situaciones lo suficientemente amplio. Intervino también un número considerable de correctores (10 en Matemáticas I y 20 en Filosofía, que estaban adscritos a los 18 tribunales de Barcelona, en junio de 1995). Esto garantizó una cierta representatividad de la muestra respecto a la población de correctores.
- Los correctores desconocían el uso que se haría de la puntuación (según la información que se dio a los correctores, ellos no podrían saber si su nota sería la oficial o si sería

utilizada solamente con finalidades estadísticas).

En los Gráficos IV y V que a continuación se recogen, se representa la puntuación que obtuvo cada examen según el corrector oficial y según el segundo corrector (NOTPAAU1 y NOTPAAU2, respectivamente). En los gráficos se informa (a través de los símbolos (o, +) de la modalidad (opción A o B) que fue escogida por el alumno. Dichos gráficos muestran una discrepancia considerable entre correctores, que por otra parte es mayor en Filosofía que en Matemáticas I. Al distinguir por opciones, parece ser que en Filosofía existe una mayor concentración de puntos de la opción A en el extremo inferior izquierdo, mientras que la opción B predomina en el extremo superior derecho. En general, los exámenes de la opción B de Filosofía han obtenido, tanto en la primera como en la segunda corrección, notas superiores a los de la opción A. En Matemáticas I no se aprecia, a primera vista, una tendencia tan clara como en Filosofía.

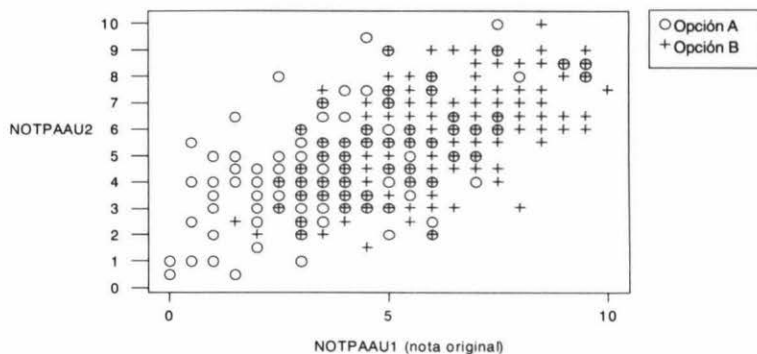
Por su parte, la Tabla II nos ofrece, para cada asignatura por separado, la correlación entre la corrección oficial, la segunda corrección y la correspondiente puntuación que en COU obtuvo cada alumno, lo que confirma la dispersión que se apreciaba en los Gráficos IV y V. Es natural que la segunda corrección mantenga una menor correlación con la nota obtenida en COU que la corrección oficial. Se trata de un efecto debido al diseño: el número de correctores es mayor en la segunda corrección que en la primera.

(31) Se fotocopiaron todos los exámenes de estas asignaturas de dos de los tribunales. Se repartieron las fotocopias aleatoriamente entre el resto de correctores. Los correctores de las fotocopias recibieron en el sobre que contenía los exámenes de su tribunal veinte fotocopias, aproximadamente, de exámenes sin corregir junto con una carta en que se les pedía que corrigieran estos veinte exámenes con los mismos criterios y al mismo tiempo que el resto de exámenes.

GRÁFICOS IV Y V

Diagramas de la doble corrección. NOTPAAU1 es la puntuación que dio el corrector oficial y NOTPAAU2 la que dio el segundo corrector

DOBLE CORRECCIÓN EN FILOSOFÍA



DOBLE CORRECCIÓN EN MATEMÁTICAS

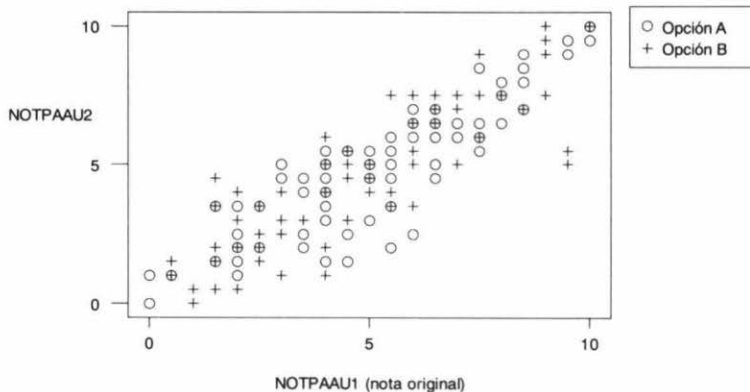


TABLA II

Correlaciones Pearson entre las notas de Filosofía y entre las notas de Matemáticas

Filosofía			Matemáticas		
	COU	NOTPAAU1		COU	NOTPAAU1
NOTPAAU1	0.316		NOTPAAU1	0.614	
NOTPAAU2	0.311	0.600	NOTPAAU2	0.572	0.074

Según la Tabla III que a continuación aparece, el grado de concordancia entre correcciones no es muy alto ³². Como era de esperar, se observa una mayor coincidencia entre correctores en Matemáticas que en Filosofía. No deja de sorprendernos, sin embargo, que en el 28% de los casos la diferencia entre las dos correcciones de Matemáticas supera la unidad. Otro dato preocupante es el siguiente: la existencia de un 13% de casos en Filosofía con una discrepancia superior a los tres puntos.

Antes hemos señalado que algunos estudios ya publicados en los que se planteaba la comparación entre la corrección de dos tribunales se habían basado en el análisis de la varianza. Por nuestra parte, hemos realizado tales análisis con nuestros datos y podemos decir que, al igual que en el estudio de Escudero y Bueno (1994) ³³, no ha resultado ser significativa la diferencia entre la primera y la segunda corrección, tanto en Filosofía como en Matemáticas. Ahora bien, aunque, en promedio la diferencia entre correcciones no sea estadísticamente significativa, puede muy bien ocurrir que, para un número importante de estudiantes, el hecho de tener

un corrector u otro modifique sus posibilidades futuras (sabemos que el acceso a algunos estudios universitarios depende de décimas). Naturalmente, hay que empezar por saber cuál es la variabilidad en que nos movemos y a partir de ahí definir el plan de trabajo a seguir: especificación de objetivos, oportuno seguimiento, control e intervención.

Si antes decíamos que al considerar la primera y la segunda corrección no se podía concluir que las medias fueran significativamente diferentes, al introducir el factor opción de examen este aspecto cambia. La opción de examen es un factor diferenciador en el siguiente sentido: la media de notas que obtienen los estudiantes de Filosofía que escogen la opción A es significativamente diferente de la media de notas que obtienen aquellos que escogen la opción B. Este hecho se da tanto en la corrección original como en la segunda corrección y en las dos materias, aunque con un grado de significación más alto en la materia de Filosofía que en la de Matemáticas. Estos análisis corroboran lo que desde el punto de vista descriptivo apuntaban los gráficos de doble corrección.

TABLA III

Frecuencias de valores de la variable diferencia (en valor absoluto) entre las dos correcciones. Para 185 estudiantes la nota de Filosofía otorgada por el primer corrector difiere en menos de un punto de la nota que le habría otorgado el segundo

	$\text{Dif} \leq 1$	$1 < \text{Dif} \leq 3$	$3 < \text{Dif}$	Total de estudiantes
Filosofía	185 (51.0%)	130 (35.8%)	48 (13.2%)	363
Matemáticas	134 (71.7%)	49 (26.2%)	4 (2.1%)	187

(32) Cabe destacar que en aquellos momentos, junio de 1995, ya se había realizado un esfuerzo considerable para adecuar los programas de COU y modificar los formatos de examen y criterios de corrección en aras de una mayor objetividad.

(33) Cabe recordar que en el trabajo de ESCUDERO y BUENO (1994) la nota de Matemáticas no se estudiaba propiamente ya que aparecía mezclada con el resto de primeras obligatorias de opción. Sí se estudiaba la nota de Filosofía.

DESCOMPOSICIÓN DE LA VARIACIÓN OBSERVADA

Con el objetivo de profundizar un poco más en el estudio de la discrepancia observada, decidimos *modelizar* nuestros datos descomponiendo el error de medida introducido en la corrección en sus diferentes fuentes de variación. Nuestro enfoque se enmarca en la teoría de la generalizabilidad (Cronbach, 1972) y en la adaptación que Longford (1994 y 1995, cap. 2) propuso para el estudio de datos relativos a exámenes y a correctores. Al intentar explicar el hecho de que dos correctores asignen puntuaciones diferentes al mismo examen podríamos distinguir entre dos posibles fuentes de discrepancia (Longford, 1994): la severidad y la inconsistencia.

Por *severidad* de un corrector, entenderemos la diferencia entre dos cantidades no observables: la media del corrector (que conoceríamos si dicho corrector corrigiera todos los exámenes) y la media global (calculable si todos los exámenes fueran corregidos por todos los correctores). No obstante, parece evidente que la discrepancia no se debe solamente a los diferentes grados de severidad. Un mismo examen puede obtener una puntuación diferente si se trata de uno de los primeros exámenes que tiene que evaluar un corrector o si por el contrario, éste se enfrenta a él cuando ya lleva corregidos un buen número de ellos. El cansancio puede influir

en la agudeza y en la atención que se pone en la corrección. También el hecho de haber visto el contenido de muchos exámenes puede modificar³⁴ el criterio de corrección que, a partir de un momento del proceso, puede volverse más indulgente o más exigente que al principio. A esta segunda fuente de error, que engloba una serie de imperfecciones que están presentes en el proceso de corrección, la llamaremos inconsistencia o *error no sistemático*. La inconsistencia específica de cada examen y corrector sería la «desviación de la puntuación otorgada respecto a la puntuación que en promedio dicho corrector otorgaría al examen en cuestión». El modelo concreto de componentes de la varianza que proponemos para explicar la variación de la puntuación de un examen es el modelo aditivo (Longford 1994):

$$y_{ij} = \alpha_i + \beta_j + \varepsilon_{ij} \quad (\text{MOV})$$

siendo $y = 1, 2, \dots, I$, el índice del examen o estudiante; $j = 1, 2, \dots, J$ el del corrector. El número de puntuaciones que entran en el estudio es $2I$; y_{ij} es la puntuación que el corrector j ha dado al examen i ; α_i es la puntuación verdadera y no observable del examen i ; β_j es la *severidad* del corrector j ; ε_{ij} representa la *inconsistencia* específica de cada corrección. Suponemos que estos tres últimos términos están mutuamente no correlacionados con medias iguales a μ , 0 y 0

(34) Uno de los hechos observados es el de la adaptación del corrector al grupo de exámenes. Parece ser que algunos profesores distribuyen a sus alumnos como harían en su propia aula o grupo-clase sin tener en cuenta que deben aplicar unos criterios universales y prescindir de un particular grupo de estudiantes que están corrigiendo. Este fenómeno de adaptación genera injusticias. En función del conjunto de centros que van a parar a un mismo tribunal puede repercutir una ventaja o inconveniente para cada alumno en particular.

Este fenómeno se limitaría si: 1) los correctores no fueran adscritos a tribunales, separando vigilancia de corrección; 2) cada corrector recibiera un bloque aleatorizado de exámenes, con desconocimiento de las escuelas de procedencia; 3) se repartieran normas consensuadas de corrección de cada examen.

Sobre las normas de corrección cabría distinguir entre las generales, elaboradas a priori y aplicables a cualquier examen, y las específicas, elaboradas por el equipo que propone los enunciados de examen en el momento de su confección y revisadas, por este mismo equipo, a partir de la corrección de una muestra aleatoria de los exámenes una vez que se dispone de los mismos.

y varianzas σ_a^2 , σ_b^2 , σ_e^2 , respectivamente. α_i sería la media que obtendríamos si todos los correctores corrigieran el examen i , mientras que β_j sería la diferencia entre la media global μ (todos los exámenes corregidos por todos los correctores) y la media correspondiente al corrector j (todos los exámenes y este corrector); ϵ_{ij} recogería la separación del corrector j en el examen i respecto de su comportamiento medio.

No hay que perder de vista que una buena corrección requiere que las componentes de la varianza relativas a la severidad y a la inconsistencia sean pequeñas con relación a la varianza de la nota verdadera.

ESTIMACIÓN

El método que hemos seguido para estimar el modelo del apartado anterior es el de los momentos, una adaptación de los clásicos estimadores de la varianza (ADEVA) al caso de datos no balanceados y para más de

un factor. Una vez determinadas las fórmulas algebraicas³⁵ para los estimadores, utilizamos el software estadístico MLn (Goldstein, 1986a) para la confección de los programas que calculan las estimaciones y permiten la revisión de las hipótesis del modelo.

En una primera comparación de las varianzas de las dos asignaturas, constatamos que en Filosofía se da una concentración de notas alrededor de su media mucho mayor que en Matemáticas. Este hecho se da tanto en la varianza total como en la varianza debida a la nota verdadera. La prueba de Filosofía en las PAAU discrimina menos que la de Matemáticas I. En cambio, la varianza de estas dos asignaturas en COU es muy similar (1.46 en Matemáticas frente a 1.36 en Filosofía)³⁶.

RESULTADOS

Éstas pueden ser algunas de las conclusiones a las que hemos llegado en nuestro trabajo:

TABLA IV
*Estimaciones de las componentes de la varianza de la puntuación observada.
Entre paréntesis aparece la proporción de varianza respecto de la varianza total*

Componentes de la varianza total	Matemáticas	Filosofía
$\hat{\sigma}_a^2$, var. entre notas verdaderas σ_i	5.350 (86.5%)	2.475 (60.2%)
$\hat{\sigma}_b^2$, var. de la severidad β_j	0.011 (0.2%)	0.248 (6 .0%)
$\hat{\sigma}_e^2$, var. de la <i>inconsistencia</i> ϵ_{ij}	0.827 (13.3%)	1.386 (33.7%)
Varianza total (suma) estimada	6.188	4.109
<i>Varianza muestral</i>	<i>6.189</i>	<i>4.065</i>

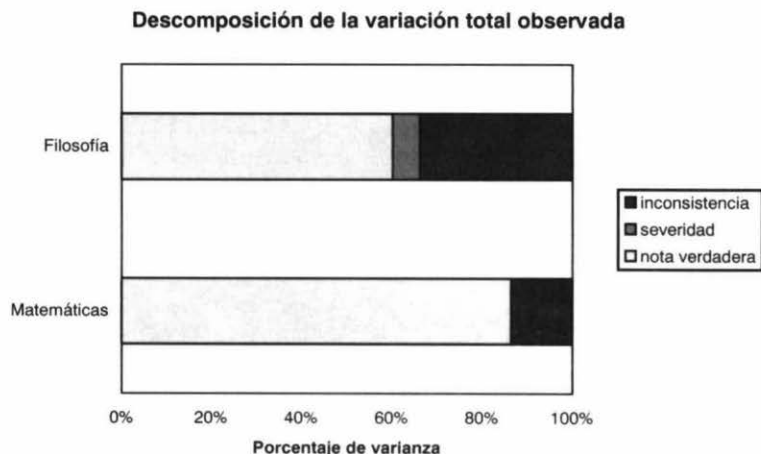
(35) Los detalles técnicos se encuentran en la tesis que A. Cuxarts está elaborando.

(36) No podemos olvidar que las puntuaciones en las asignaturas de COU son en cuatro categorías (5.5, 6.5, 7.5 o 9), hecho que provoca una excesiva similitud entre estudiantes y entre distribuciones por asignaturas en COU.

- El estudio arroja esta primera conclusión: la calidad de la corrección es baja. Para una tercera parte de los exámenes de Filosofía y un 18% de Matemáticas I, la diferencia entre las dos correcciones superó los dos puntos. Los indicadores de la calidad de la corrección, obtenidos a partir de la *modelización* de los datos –léase, las varianzas de la *nota verdadera*, de la *severidad* y de la *inconsistencia*, así como las correlaciones–, confirman las primeras observaciones derivadas de la simple comparación de puntuaciones.
- En las dos materias analizadas, la principal fuente de error en la corrección resulta ser la inconsistencia. La varianza debida a la inconsistencia representa un 13% de la varianza total en Matemáticas y un 34% en Filosofía. En Matemáticas no se aprecia una diferencia entre la severidad de los correctores. Sin embargo, en cambio en Filosofía parecen coexistir diferentes grados de severidad entre correctores –la varianza de la severidad en Filosofía representa un 6% de la varianza total. El Gráfico VI ilustra sobre la participación porcentual de cada fuente de error en la variación total.
- El análisis de la inconsistencia apunta la posibilidad de que ciertos exámenes o preguntas conlleven mayor probabilidad de discrepancia entre correctores. El hecho que más llama la atención, que sustenta nuestra conjetura, es la evidencia de que en Filosofía una de las opciones ha generado más discrepancia que la otra y, además, lo ha hecho en sentido inverso, de tal manera que, la diferencia entre las puntuaciones del primer y segundo corrector en promedio es más del doble en la opción A (-0.714) que en la B (0.33). Los diagramas de caja (box-plot) que se recogen en el Gráfico VII ilustran esta situación. En ellos, la línea que divide cada caja se sitúa en la diferencia mediana. El diagrama relativo a la opción B se desplaza hacia valores más positivos,

GRÁFICO VI

Descomposición de la varianza total de la nota observada en Filosofía y en Matemáticas I, según aplicación del modelo (MDV)



lo cual pone de manifiesto que para esta opción la corrección oficial tendió a ser superior a la de la segunda corrección. Por el contrario, para los exámenes de la opción A fue la segunda corrección la que tendió a ser superior. Una de las razones de esta situación -según se desprende de un estudio más pormenorizado sobre los exámenes de los alumnos y las anotaciones de los correctores- parece ser el comportamiento diferenciado de los dos correctores oficiales que participaron en el estudio. Por ello, se puede deducir que parte de la diferencia observada entre tribunales se debe, pues, al comportamiento diferenciado de sus respectivos correctores (oficiales).

- En cuanto a los efectos que las correcciones hubieran tenido en el acceso a la universidad, se puede decir que aproximadamente un 3% de los estudiantes de la muestra (11 sobre un total de 362 alumnos) habrían sido ubicados de manera diferente (obtener o no un aprobado en la nota de acceso) si sus exámenes hubieran sido

evaluados por los segundos correctores en lugar de por los correctores oficiales.

En condiciones normales solamente se dispone de una corrección por examen. Es de destacar la mejora que se introduciría, sobretodo en Filosofía, si para cada examen se pudiera contar con dos correcciones y tomar como nota definitiva la media de ambas. Según se desprende de la Tabla V, el valor del coeficiente de correlación entre la puntuación observada y la nota verdadera, se vería incrementado en un 12% al tomar la media entre las dos correcciones (de 0.78 a 0.87). En Matemáticas tan sólo representaría una mejora del 3%.

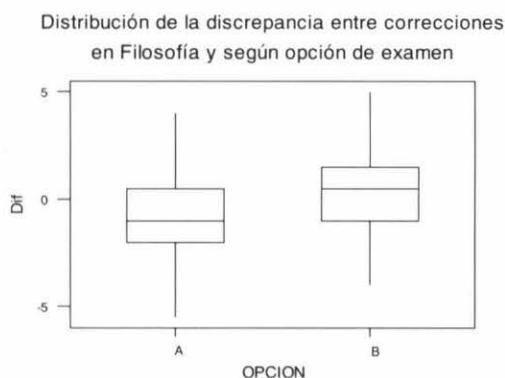
Por último hay que hacer constar que, a pesar que los resultados de este estudio no contradicen la intuición, no pueden generalizarse a otras materias, ni a anteriores o futuras convocatorias, ni tampoco a diferentes formatos de examen.

ALGUNAS REFLEXIONES Y PROPUESTAS

El Consejo de Universidades en un documento del año 1993 que lleva por título

GRÁFICO VII

*Discrepancia entre correcciones en Filosofía y opción de examen
(DIF = Corrección oficial - segunda corrección)*



«Las Pruebas de Aptitud para el acceso a la Universidad: Problemática actual y propuesta de modificación», señala concretamente la necesidad de una doble corrección para la prueba del Comentario de Texto (propuesta 3.3):

Dadas las características especiales de este examen, sería conveniente utilizar para el Comentario de Texto una doble corrección. La doble corrección es utilizada en distintos tipos de prueba y ejercicios como garantía de ponderación y equilibrio en la calificación final que se otorga.

En el caso de la prueba de Comentario de Texto, que no se refiere a una materia específica, sino que tiene un carácter general, y dado que los profesores que han de corregirlo, al ser de cualesquiera de las otras materias que integran la prueba, pueden tener apreciaciones y criterios dispares, se hace más aconsejable el uso de este procedimiento.

La calificación final será la media de la otorgada por ambos correctores, siempre que no se diferencien en más de 2 puntos. Si la diferencia es mayor, los correctores procederían, en cada caso, a la revisión de las calificaciones efectuadas siguiendo meticulosamente los criterios de corrección específicos establecidos para el examen. La corrección de este examen se dará por concluida una vez que todas las calificaciones se hayan obtenido promediando puntua-

ciones con diferencias que no superen el límite establecido (Modificación legal del R.D. 406/83. /Art. 4.1. (nuevo) y de la Orden 3-9-37 /Art.nuevo).

A este respecto, y teniendo en cuenta todos los aspectos que aquí hemos apuntado cabe plantearse, a nuestro entender, tres observaciones relevantes:

- Esta propuesta no se está teniendo en cuenta en las actuales Pruebas de Aptitud o Acceso a la Universidad, al menos en todas las administraciones con las que hemos contactado.
- Si para una materia como el Comentario de Texto se considera no admisible una diferencia superior a 2 puntos entre las dos correcciones, podemos preguntarnos cuál habría de ser la diferencia que se podría admitir entre las dos correcciones en materias como Filosofía o Matemáticas, donde los correctores son expertos en el tema.
- El presente estudio, realizado sobre las pruebas de acceso en las materias de Matemáticas y de Filosofía ha arrojado los siguientes datos: ³⁷ 14 exámenes de Matemáticas (un 7% de los mismos) recibieron puntuaciones que diferían en más de 2 puntos, mientras que en la materia Filosofía

TABLA V

Indicadores de la calidad de la corrección, calculados a partir de las estimaciones que se derivan de la aplicación del modelo (MDV)

Coeficientes de correlación	Matemáticas	Filosofía
r (entre las dos puntuaciones)	0.86	0.60
$r_{\alpha 1}$ (entre <i>nota verdadera</i> y puntuación observada)	0.93	0.78
$r_{\alpha 2}$ (entre <i>nota verdadera</i> y puntuación media)	0.96	0.87

(37) La distribución de estos exámenes por opciones A / B fue de 10/4 de Matemáticas y 33/44 en Filosofía, respectivamente.

fueron 77 los exámenes cuya calificación distaba en más de dos puntos en cada una de las correcciones (lo que supone un 41%).

A la vista de estos datos y con el objetivo de incrementar la precisión en la corrección de las pruebas PAAU, presentamos, a continuación, una serie de propuestas:

- Considerar en cada materia de examen la posibilidad de sustituir el examen actual, o una parte del mismo, por una prueba con preguntas de respuesta cerrada.
- De manera sistemática, y para los exámenes con preguntas de repuesta abierta, en cada convocatoria se debería seleccionar una muestra de exámenes de cada materia y realizar una doble corrección de los mismos con el objetivo de medir la fiabilidad o precisión en su corrección y detectar posibles fuentes de discrepancia.
- Para aquellas asignaturas con una precisión baja realizar una doble corrección de todos los exámenes, como ya recomendaba el Consejo de Universidades en 1993.
- Incorporar un mecanismo de revisión automático de todas las puntuaciones PAAU que, al comparar con las de la misma prueba en COU u otra variable indicativa, permita destacar aquellas puntuaciones que se separan demasiado de las previsiones. Realizar una doble corrección (si se trata de preguntas de repuesta abierta) de estos casos e introducir los ajustes que se consideren oportunos. Aquellos centros (y correctores) cuyos alumnos (exámenes) en un porcentaje alto hayan sufrido revisión deberían ser analizados.

No obstante, creemos oportuno traer a consideración el hecho de que para poder llevar a cabo de manera eficaz todas las propuestas anteriores habría que tener una infraestructura y unos medios adecuados.

CONCLUSIONES

Destacamos a continuación las principales conclusiones que se derivan de la aplicación de modelos de regresión con coeficientes aleatorios para el estudio de la asociación entre la nota de COU y la nota de las PAAU de cada estudiante:

- Existe una variación significativa de la nota PAAU entre centros escolares. Un 20% aproximadamente de la variación total de la nota PAAU corresponde a variación entre centros.
- La influencia del centro escolar en la predicción de la nota PAAU individual se concreta en un término aditivo, común a todos los estudiantes del mismo centro y que hemos llamado efecto centro. Tales efectos tienen asociada una distribución de probabilidad y permiten identificar los centros escolares que presentan una asociación entre la nota COU y la nota PAAU extrema.
- Las distribuciones de la nota COU, a diferencia de la nota PAAU, varían muy poco de un centro a otro. De ahí que el coeficiente de correlación intra-centros para la nota COU sea prácticamente 0. Este hecho sugiere que los centros están utilizando escalas de puntuación propias, diferentes de un centro a otro, según pone en evidencia el examen PAAU.
- El modelo de regresión de coeficientes aleatorios de la nota PAAU frente a la nota COU que contempla género, posible repetición de COU, opción de COU y tipo de centro, nos lleva a una serie de conclusiones en cuanto al valor predictivo de estas variables coincidentes con anteriores estudios realizados en el ámbito estatal. La novedad de nuestro enfoque estriba en la determinación del papel predictivo de cada

centro en la nota individual de PAAU de manera conjunta con el resto de variables citadas.

Por otra parte, el análisis del proceso de corrección de las pruebas PAAU ha puesto en evidencia los siguientes aspectos:

- La baja calidad de la corrección en las dos asignaturas estudiadas, incluso en Matemáticas. La principal fuente de error ha sido la *inconsistencia* (la varianza debida a la *inconsistencia* representa un 13% de la varianza total en Matemáticas y un 34% en Filosofía). En esta última materia, además parecen coexistir diferentes grados de *severidad* entre correctores (un 6% de la varianza total corresponde a *severidad*).
- Son importantes las consecuencias que para algunos estudiantes se pueden derivar de esta imperfección del sistema: un 3% de los estudiantes de la muestra habrían tenido una suerte distinta en su incorporación al mundo de la universidad de haberse tenido en cuenta la segunda corrección en lugar de la corrección oficial. Los estudiantes más afectados por la baja fiabilidad son, naturalmente, los que se encuentran cerca de la frontera (*borderline*) del aprobado.
- El valor que tiene el monitorizar una investigación conectada a la ejecución,³⁸ tanto por la información que suministra como por la posibilidad de intervenir para realizar un ajuste a tiempo.
- La necesidad de interpenetrar correctores y exámenes si queremos comparar los resultados de diferen-

tes tribunales, centros, comarcas, etc.

- La conveniencia de contrastar empíricamente la dificultad de las preguntas de cada examen y materia.
- La existencia de diversas fuentes de variación en la corrección: algunas de ellas relacionadas con el diseño y el contenido de los exámenes (de ahí la necesidad de mejorar el procedimiento de elaboración de los exámenes); otras relacionadas con la organización de las pruebas (por este motivo recomendamos la separación entre la labor de vigilante y la labor de corrector).
- La importancia que adquieren todos estos temas en la discusión de las nuevas PAAU. En la actualidad, el hecho de utilizar una puntuación que es una «media de medias» diluye, en gran parte, los efectos de un sistema imperfecto. Una ponderación que diera un peso mayor a algunas de las materias requeriría un mayor grado de fiabilidad en la corrección de las mismas y justificaría la incorporación de pruebas de respuesta cerrada.

PERSPECTIVAS DEL TRABAJO DE INVESTIGACIÓN

Los modelos de descomposición de la varianza han demostrado ser de utilidad en la investigación realizada hasta el momento y consideramos interesante ahondar en sus posibilidades.

En una de las investigaciones que ahora se está realizando hemos abordado el estudio de la variación conjunta del vector de notas PAAU de cada estudiante. Con

(38) En (CUXART and LONGFORD, 1996) se incluyen una serie de reflexiones y propuestas sobre el efecto de la elección, la posibilidad de realizar *reajustes*, de comparar resultados, y de realizar *pretests* de las preguntas para conocer su dificultad.

ello tratamos de detectar la asociación existente entre las diversas materias y el poder discriminador del primer y del segundo ejercicio y de la nota global de PAAU. También intentamos obtener información sobre la capacidad evaluadora de ambos ejercicios.

Los resultados aportados han sido éstos: las correlaciones entre las diferentes pruebas son muy débiles, incluso cuando se calculan para los estudiantes de una misma opción. Es un hecho conocido que el error de medida en cada evaluación tiene el efecto de atenuar los coeficientes que miden la relación entre variables y que las preguntas de respuesta abierta conllevan mayor subjetividad e imprecisión en su corrección. Si además, los formatos de examen de dos asignaturas, aunque sean propias de la opción, son muy diferentes pueden estar evaluando habilidades distintas a la vez que conocimientos. Por otro lado, los actuales exámenes no cubren de manera exhaustiva el programa de las asignaturas. De ahí que pueda hablarse de un factor suerte en cuanto a los temas que aparecen cada año a examen. La suerte de una asignatura a otra puede variar y nos encontramos con otra fuente de variabilidad. Todas estas consideraciones hacen referencia a la validez del examen y a la fiabilidad del mismo.

En el futuro creemos que sería interesante la utilización de modelos estadísticos que tengan en cuenta el error de medida. Para ello necesitamos disponer de réplicas (doble corrección, por ejemplo) para al menos una muestra de cada materia. Pensamos que los modelos LISREL nos permitirían introducir un poco más de luz en el complejo sistema de relaciones entre las materias y entre los factores que influyen en su evaluación.

Del análisis de la corrección que se ha hecho teniendo en cuenta las dos opciones

(A o B) de las materias de Filosofía y Matemáticas I se desprende que el nivel de puntuaciones no es el mismo en las dos opciones. En general, y tanto en la primera como en la segunda corrección, la opción A de Filosofía recibió notas inferiores a la opción B, mientras que en Matemáticas ocurrió la situación inversa. Dado que los estudiantes fueron quienes eligieron la opción de examen es imposible separar el factor opción de la preparación del estudiante. Se plantea pues la necesidad de conocer la dificultad ³⁹ de las preguntas planteadas y, a la vez, recomendar la limitación al máximo de la opcionalidad en estos exámenes.

La implantación de la LOGSE y las nuevas PAAU representa una oportunidad para introducir cambios estructurales en el proceso de acceso a la universidad que tengan en cuenta la experiencia acumulada con el sistema de acceso anterior. Consideramos un hecho clave la incorporación de la investigación estadística en el seguimiento del proceso y la evaluación empírica de las modificaciones que se vayan introduciendo.

En una situación de transición como la que estamos viviendo en estos momentos en la que existen dos sistemas de pruebas PAAU vigentes, es interesante conocer y comparar porcentajes de superación de las diferentes fases educativas en relación a la población de jóvenes de una misma edad o cohorte. De ahí también surge la oportunidad de investigar la aplicación de modelos de regresión logística.

ALGUNAS CONSIDERACIONES PEDAGÓGICAS A LA LUZ DE LAS ESTADÍSTICAS

Entre las consideraciones pedagógicas que se desprenden del presente estudio,

(39) Este curso se están analizando los datos de un nuevo experimento de doble corrección para un número mayor de asignaturas. El objetivo de este experimento es calcular la fiabilidad de la corrección, conocer la dificultad de las preguntas y recoger la opinión de los correctores en cuanto al enunciado de examen y a los criterios de corrección específicos.

consideramos oportuno destacar las siguientes:

- La relevancia que adquiere para todo sistema educativo el hecho de disponer de datos fiables. En este sentido las pruebas PAAU, como examen externo a los centros y estándar,⁴⁰ se revelan como un instrumento de gran utilidad.
- La necesidad de comparar resultados con rigor y teniendo en cuenta el contexto. Si se dan estos dos requisitos, la prevención a la comparación entre centros carecerá de sentido.
- Se debería avanzar en la «cultura» de realizar estudios que sean útiles para la Administración y que, al mismo tiempo, sirvan de referencia y contraste para los centros.
- La defendida necesidad de proponer exámenes que contengan preguntas lo más cerradas que se pueda plantea un cambio en la pedagogía. Los profesores de secundaria deberían incorporar en su docencia este tipo de pruebas.⁴¹
- Las pruebas PAAU no debe ser el primer examen global de la materia al que se enfrentan los estudiantes. Por ello, es importante que en la enseñanza secundaria los alumnos preparen y realicen exámenes que abarquen una parte importante de la programación, y a ser posible, todos los contenidos del temario.
- Las habilidades en comunicación escrita adquieren una gran relevancia en las pruebas de respuesta abierta de las PAAU —hemos hablado de limitar en las PAAU las pruebas de respuesta abierta, no de

eliminarlas—. Y puesto que dichas habilidades requieren un proceso de aprendizaje, resulta evidente que esta necesidad debería ser atendida como una de las prácticas prioritarias en la enseñanza secundaria.

BIBLIOGRAFÍA

- AGUIRRE DE CÁRCER: *La Selectividad a debate*, Madrid, Universidad Autónoma, 1984.
- AITKIN, M. y LONGFORD, N.: «Statistical modelling issues in school effectiveness studies», en *J. R. Statistical Society*, 149 (1986), Part 1, pp. 1-43.
- COCHRAN, W. G.: *Sampling Techniques*, 3.^a ed. Toronto, Wiley, 1977.
- CRONBACH, L. J.; GLESER, G. C.; NANDA, H. y RAJARATNAM, N.: *The dependability of behavioral measurements: Theory of generalizability for scores and profiles*, New York, John Wiley, 1972.
- CUXART, A.; GRAFFELMAN, J. y MARTÍ, M.: *La nota PAAU y su relación con la nota COU: un modelo de regresión con coeficientes aleatorios para el estudio del efecto centro en la nota PAAU*, Actas de la 5.^a Conferencia Española de Biometría. Valencia. 1995.
- CUXART I JARDÍ, A. y LONGFORD, N. T.: *Equity in the university admissions process in Spain*, Estudio —en prensa— presentado al 10th European Congress of Psychometric Society, Santiago de Compostela, julio de 1997 y al 5th European Research Congress on Education, Frankfurt, Alemania, set. 1997.
- ESCUADERO, T.: «Investigaciones y Experiencias: Buscando una mejor selección de

(40) En el artículo «El desarrollo de la LOGSE: las nuevas pruebas de acceso a la universidad» de este mismo número se dan diferentes argumentos a favor de las pruebas externas.

(41) Existen múltiples experiencias en este ámbito: Las pruebas *Canguro* de matemáticas que se realizan simultáneamente en varios países europeos, las pruebas MIR de nuestro país, el examen SAT americano,...

- universitarios», *Revista de Educación*, núm. 283, Ministerio de Educación y Ciencia, 1987.
- ESCUADERO, T. y BUENO GARCÍA, C.: «Investigaciones y Experiencias: Examen de Selectividad. El estudio del tribunal paralelo», *Revista de Educación*, núm. 304, Ministerio de Educación y Ciencia, 1994.
- GOLDSTEIN, H.: «Multilevel mixed linear model analysis using iterative generalized least squares», *Biometrika*, 73, pp. 43-56, 1986a.
- *Multilevel models in Educational and Social Research*, Oxford University Press, New York, 1987.
- GOLDSTEIN, H. y SPIEGELHALTER, D. J.: «League tables and their limitations: statistical issues in comparisons of institutional performances (with discussion)», *Journal of the Royal Statistical Society*, Ser. A, 159, 1989, pp. 385-443.
- GOLDSTEIN, H.: *Multilevel Statistical Models*, 2.^a ed., Kendall's Library of Statistics 3 (London, Edward Arnold), 1995.
- LONGFORD, N. T.: *Random Coefficient Models*, Oxford Science Publications, Clarendon Press, Oxford, 1993.
- «Reliability of essay rating and score adjustment», *Journal of Educational and Behavioral Statistics*, Vol 19, núm. 3, 1994, pp. 171-200.
- «Models for uncertainty in Educational Testing», *Springer Series in Statistics*, New York, 1995.
- LÓPEZ, M.^a del R.: «Algunos resultados sobre las Pruebas de Acceso a la Universidad Autónoma de Madrid», en ÁLVAREZ, J. B. y ARROYO, F. (comp): *Acceso a la Universidad y Marco Educativo*, Tarbiya, número extraordinario, junio, 1977.
- MARTÍ RECOBER, M.: *Los sistemas de corrección de las pruebas de Selectividad en España. Análisis y propuestas. Concurso nacional de Proyectos de Investigación Educativa*, Ministerio de Educación y Ciencia, CIDE, 1995.
- MEMORIA DE ACTIVIDADES DEL CONSEJO DE UNIVERSIDADES. Junio 1991- Julio 1993.
- MUÑOZ-REPISO IZAGUIRRE, M., et al: *Las calificaciones en las pruebas de aptitud para el acceso a la universidad*, núm. 61 colección INVESTIGACIÓN. Madrid: CIDE, 1991.
- MUÑOZ-REPISO, M. y otros: *El sistema de acceso a la Universidad en España: tres estudios para aclarar el debate*. Madrid: CIDE, 1997.
- SANS, A.: *Selectivitat universitària. Anàlisi a Catalunya. Tesi doctoral*. Bellaterra: Publicacions de la UAB, 1990.
- SATORRA, A. y UDINA, F.: *Comunicación personal*, 1994.